

Comparative genomics tools for biological discovery

Inna Dubchak, Ph.D.

Staff scientist

Lawrence Berkeley National Laboratory

ildubchak@lbl.gov

Outline

What is comparative genomics?

VISTA tools developed for comparative genomics.

Related biological stories

Large scale VISTA applications including automatic
computational system for comparing the human
and mouse genomes

The Human genome



From the Nature paper:

The next steps:

Developing the IGI (integrated gene index) and IPI (integrated protein index)

Large-scale identification of regulatory regions

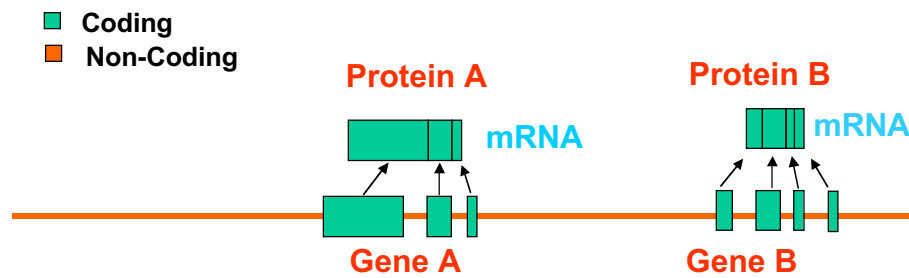
Sequencing of additional large genomes

Completing the catalogue of human variation

From sequence to function

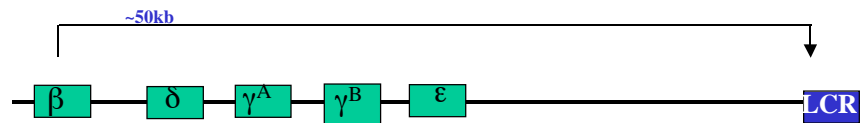
"

1-2% Coding



Distant Non-Coding Sequences Causing Disease

β-Thalassemia



Disease	Gene	Distance
Campomelic displasia	SOX9	850kb
Aniridia	PAX6	125kb
X-Linked Deafness	POU3F4	900kb
Saethre-Chotzen syndrome	TWIST	250kb
Rieger syndrome	PITX2	90kb
Split hand / split foot malformation	SHFM1	450kb

Background

Evolution can help!

In general, functionally important sequences are conserved



Conserved sequences are functionally important



Raw sequence can help in finding biological function

Comparison of 1196 orthologous genes (Makalowski et al., 1996)

- Sequence identity:
 - exons: 84.6%
 - protein: 85.4%
 - introns: 35%
 - 5' UTRs: 67%
 - 3' UTRs: 69%
- 27 proteins were 100% identical

Integrating data into more powerful gene prediction
models than with human genomic sequence alone

Comparing sequences of different organisms



- Helps in gene predictions
- Helps in understanding evolution
- Conserved between species non-coding sequences are reliable guides to regulatory elements
- Differences between evolutionary closely related sequences help to discover gene functions

Challenges

- Sequence at different stages of completion, difficult to compare

Whole genome shotgun
Finished BACs



Partial Assemblies

- Fast and accurate analysis
- Scaling up to the size of whole genomes



VISIA

USE VISIA on the WEB

 [- instructions for using VISIA](#)
 [- instructions for using VISIA](#)

DOWNLOAD VISIA

 [- download VISIA](#)
 [- download VISIA](#)

INFORMATION about VISIA

 [- VISIA](#)
 [- VISIA](#)

VISUALIZATION TOOL FOR ALIGNMENTS

VISIA is a homepage for VISIA, Visualization Tool for Alignments.

VISIA is an integrated computational system for global alignment and visualization, alignment comparison, analysis. It allows for the visualization of long sequence alignments of DNA from two of more species with annotation information, and it was developed by Odeh Costumed sequence alignment engine. [Boukhalil et al., 2007](#)

It has a rich GUI, allowing for easy distribution of sequence similarity and differences, entity level comparisons, interactive visualization of alignment of various lengths at different levels of resolution.

This system consists of several submodules:

- align** for global alignment of DNA sequences of arbitrary length. It utilizes an algorithm from dynamic programming. It can also handle one sequence in a hierarchical end-to-end oriented graph format. [2008]
- vis** is a computational tool for computing an arbitrary number of genetic

Modules of VISTA:

- 6

Sequence comparisons. How?

Three variations:

Find the best **OVERALL** alignment.

Global alignment

Find **ALL** regions of similarity.

Local alignment

Find the **BEST** region of similarity.

Optimal local alignment

Aligning large genomic regions

- Long sequences lead to memory problems
- Speed becomes an issue
- Long alignments are very sensitive to parameters
- Draft sequences present a nontrivial problem
- Accuracy is difficult to measure and to achieve

References for some existing programs:

Glass:

Domino Tiling, Gene Recognition, and Mice.

Pachter, L. *Ph.D. Thesis, MIT* (1999)

Human and Mouse Gene Structure: Comparative Analysis and Application to Exon Prediction.

Batzoglou, S., Pachter, L., Mesirov, J., Berger, B., Lander, E. *Genome Research* (2000).

MUMmer

Delcher, A.L., Kasif S., Fleischmann, R.D., Peterson J., White, O. and Salzberg, S.L.

Alignment of whole genomes. *Nucleic Acids Research* (1999)

PipMaker

PipMaker: A Web Server for Aligning Two Genomic DNA Sequences.

Scott Schwartz, Zheng Zhang, Kelly A. Frazer, Arian Smit, Cathy Riemer, John Bouck, Richard Gibbs, Ross Hardison, and Webb Miller. *Genome Research* (2000)

Scan2

Dbscan/Scan2: Fast alignment of mega-sequences.

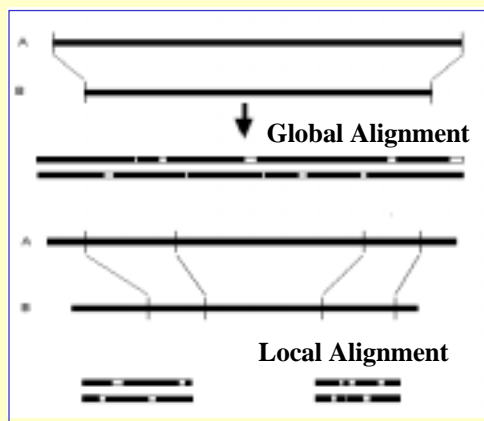
Seledtsov I.A., Solovyev V.V. To Appear. Web site <http://softberry.com/>

Local alignment algorithms are designed to search for highly similar regions in two sequences that may not be highly similar in their entirety. The algorithm works by first finding very short common segments between the input sequence and database sequences, and then expanding out the matching regions as far as possible.

For cross-species comparison one needs to accurately align two complete sequences. It is insufficient to find common similar regions in the two sequences, rather, what is needed is a global map specifying how the two sequences fit together, much like understanding how the pieces in a puzzle connect up with each other.

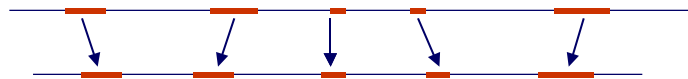
This problem is called global alignment

Local vs global alignment



AVID- the alignment engine behind VISTA

- **Very fast** global alignment of megabases of sequence.
- **Provides details** about ordered and oriented contigs, and accurate placement in the finished sequence.
- **Full integration** with repeat masking

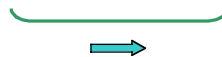


- **ORDER and ORIENT**
- **FIND** all common k-long words (k-mers)
- **ALIGN** k-mers scoring by local homology
- **FIX** k-mers with good local homology
- **RECURSE** with smaller k (shorter words)

Visualization



```
tggtacattcaaattatg-----ttctcaaagtgagcatgaca-acttttttccatgg
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
tgatgacatctatttgcgttttccttttttagaaactgcatgagagcctggcctagtaggg
```



Window of length **L** is centered at a particular nucleotide in the base sequence

Percent of identical nucleotides in **L** positions of the alignment is calculated and plotted

Move to the next nucleotide

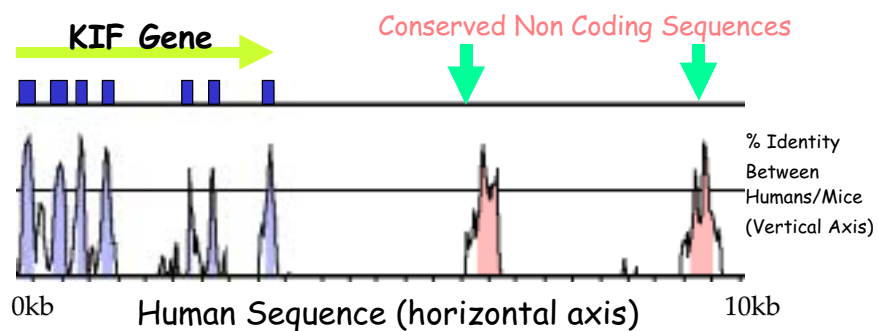
Finding conserved regions with percentage and length cutoffs

Conserved segments with percent identity X and length Y - regions in which every contiguous subsegment of length Y was at least $X\%$ identical to its paired sequence. These segments are merged to define the conserved regions.

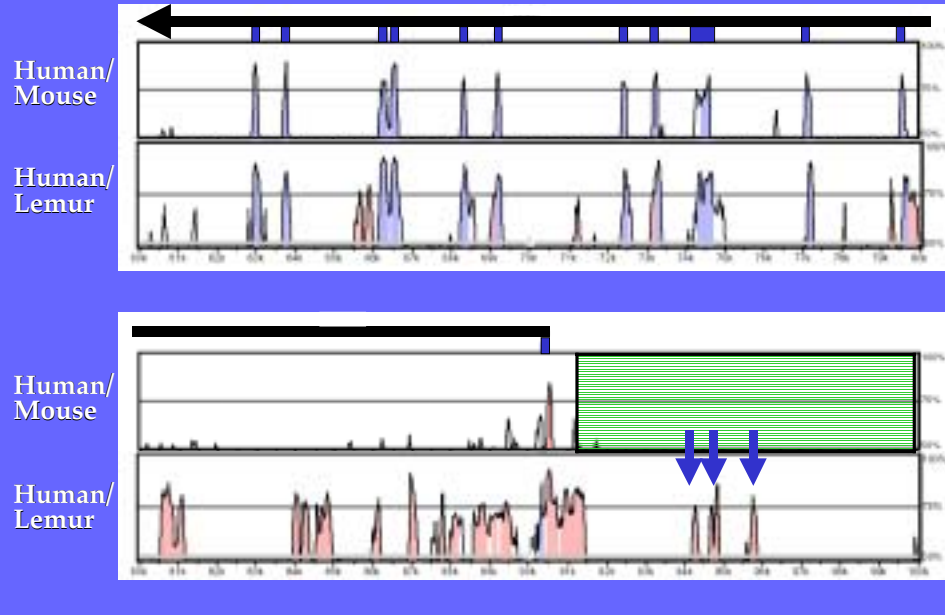
Output:

11054 - 11156 = 103bp at 77.670%	NONCODING
13241 - 13453 = 213bp at 87.793%	EXON
14698 - 14822 = 125bp at 84.800%	EXON

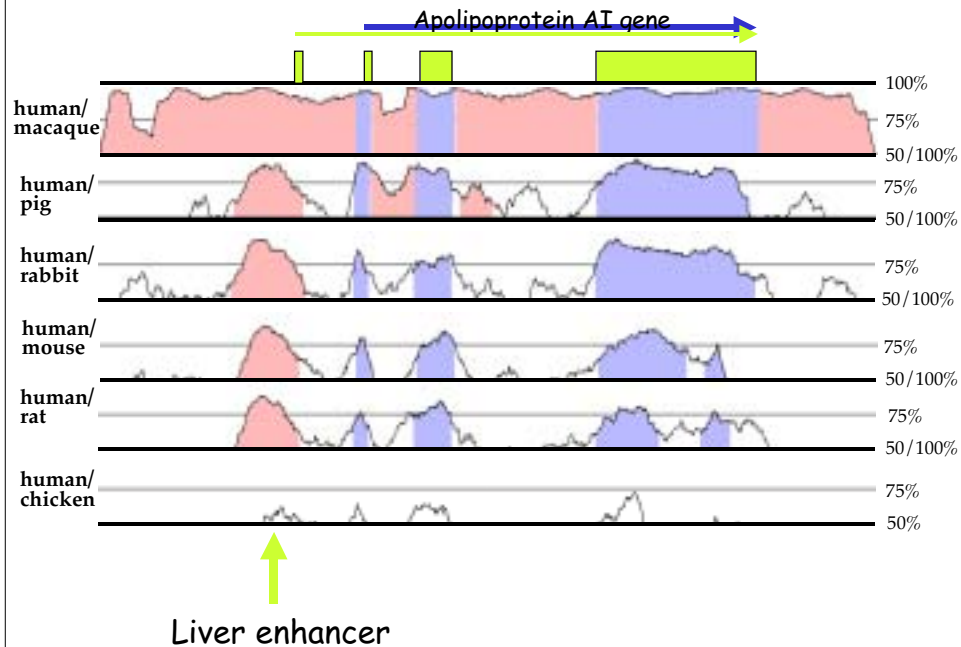
VISTA plot



Low-Density Lipoprotein Receptor (LDLR)



Multi-Species Comparative Analysis (VISTA)



VISTA input files

- Sequences of two or more organisms
- Annotation file for a base sequence if available

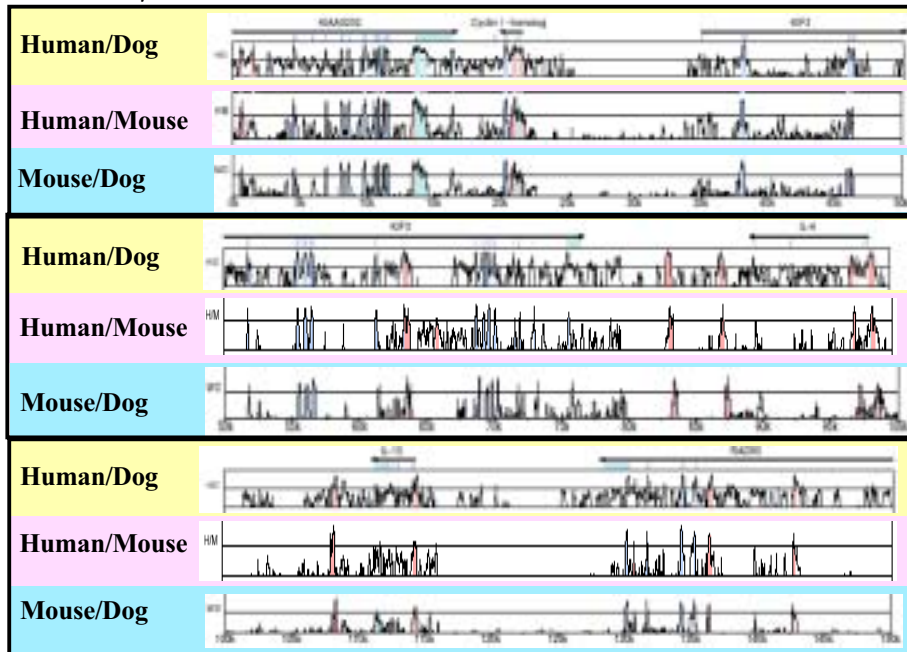
VISTA output files

- All pair wise global alignments of the sequences
- VISTA plot
- The list of conserved regions at predefined by the user length and conservation cutoffs

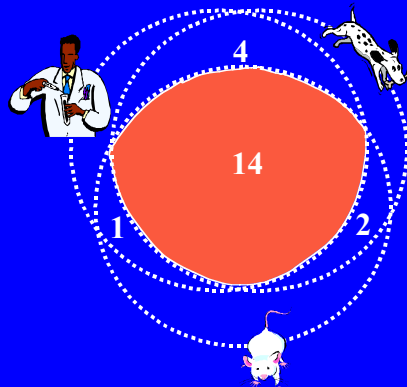
VISTA flavors

- **VISTA** - comparing DNA of multiple organisms
- **for 3 species** - analyzing cutoffs to define actively conserved non-coding sequences
- **cVISTA** - comparing two closely related species
- **rVISTA** - regulatory VISTA

Example: Dubchak et al., 2000, *Genome Research*, 10: 1304-1306.



Active conservation of noncoding sequences -
present in more than two mammals



% Cutoff

sum of three pair wise
Intersection/Union
values is maximal

Over 120 basepairs:

H/D > 92%

H/M > 80%

D/M > 77%

Identifying non-coding sequences (CNSs) involved in transcriptional regulation



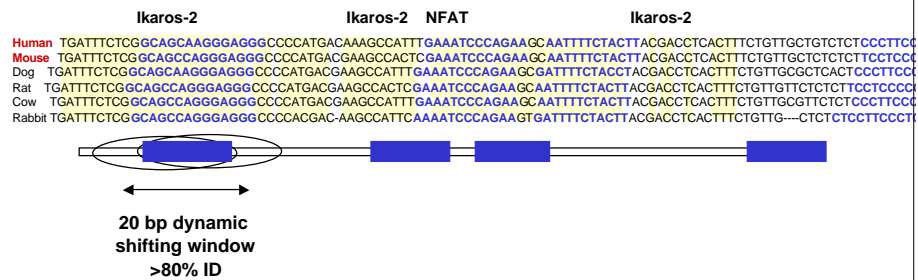
rVISTA - prediction of transcription factor binding sites

- Simultaneous searches of the major transcription factor binding site database ([Transfac](#)) and the use of global sequence alignment to sieve through the data
- Combination of [database searches](#) with [comparative sequence analysis](#) reduces the number of predicted transcription factor binding sites by several orders of magnitude

Regulatory VISTA (rVISTA)

1. Identify potential transcription factor binding sites for each sequence using library of matrices (TRANSFAC)
2. Identify aligned sites using VISTA
3. Identify conserved sites using dynamic shifting window

Percentage of conserved sites of the total 3-5%

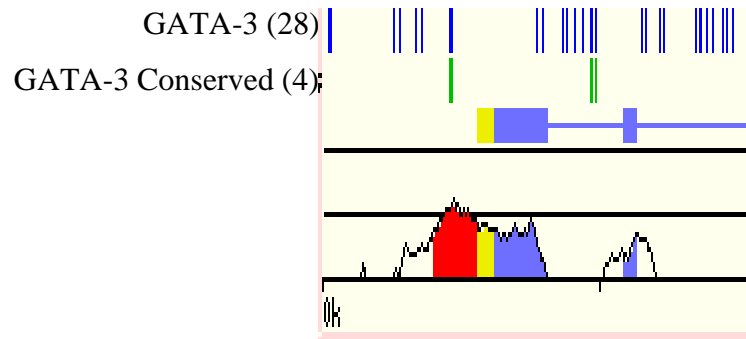


~1 Meg region, 5q31

	Coding	Noncoding
Human interval Transfac predictions for <i>GATA</i> sites	839	20654
Aligned with the same predicted site in the mouse seq.	450	2618
Aligned sites conserved at 80% / 24 bp dynamic window	303	731
Random DNA sequence of the same length		29280

2 Exp. Verified GATA-3 Sites

IL 5



Choose Families to Visualize

<input type="checkbox"/> AHR	<input type="checkbox"/> E2	<input type="checkbox"/> LIMCCOM	<input type="checkbox"/> BOB1
<input type="checkbox"/> AHRBNT	<input type="checkbox"/> E2F	<input type="checkbox"/> EYF1	<input type="checkbox"/> BOB2
<input type="checkbox"/> AML1	<input type="checkbox"/> E4	<input type="checkbox"/> MAX	<input type="checkbox"/> BOD1
<input type="checkbox"/> AP1	<input type="checkbox"/> E4BP4	<input type="checkbox"/> MEF3	<input type="checkbox"/> BDRPC4
<input type="checkbox"/> AP1K3	<input type="checkbox"/> EGR1	<input type="checkbox"/> MIF1	<input type="checkbox"/> B8
<input type="checkbox"/> AP2	<input type="checkbox"/> EGR2	<input type="checkbox"/> MYB	<input type="checkbox"/> BDP1
<input type="checkbox"/> AP4	<input type="checkbox"/> EGR3	<input type="checkbox"/> MYC/MAX	<input type="checkbox"/> BOD5
<input type="checkbox"/> ARNT	<input type="checkbox"/> ELK1	<input type="checkbox"/> MYO10	<input type="checkbox"/> BP1
<input type="checkbox"/> ARPI	<input type="checkbox"/> ER	<input type="checkbox"/> MZF1	<input type="checkbox"/> BDRP1
<input type="checkbox"/> ATF	<input type="checkbox"/> ETV1	<input type="checkbox"/> NF1	<input type="checkbox"/> BDP

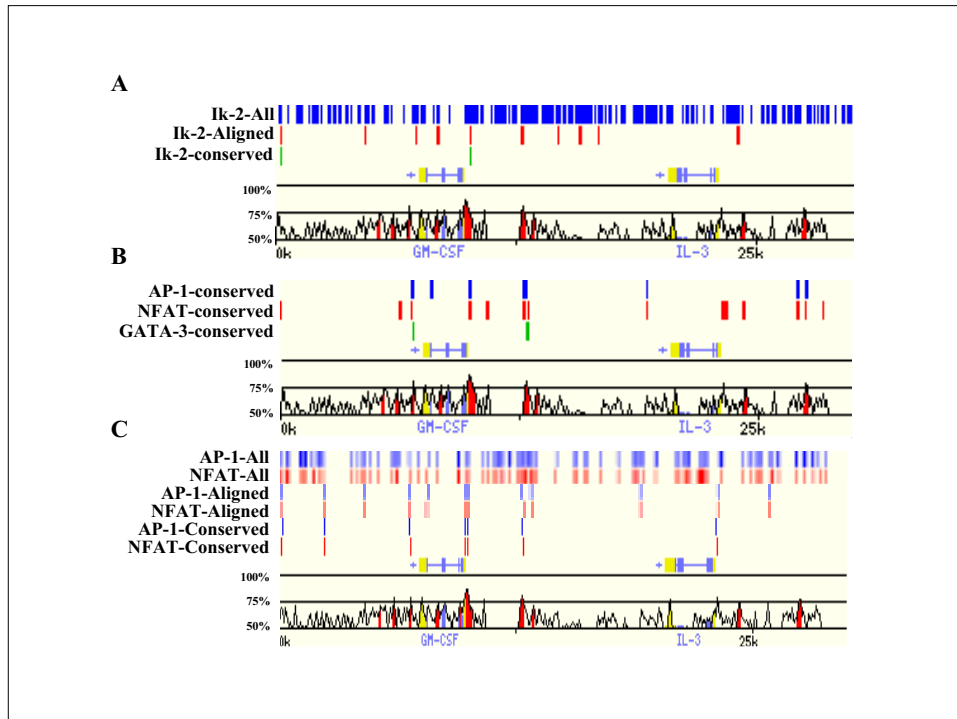
Picture
 Rows per screen: 1000
 Rows per screen per picture: 100

Clustering
 Clustering: 2 clusters per 1000 rows
 1) color per 2 clusters 2) color clustering
☐ 250 ☐ 500 ☐ 1000

Type
☐ unclustered
☐ clustered
☐ 25

Select Picture Properties, Clustering, and Type of Visualization





Sequence motif recognition

+

multiple sequence alignment of syntenic
regions,



a high throughput strategy for filtering and
prioritizing putative DNA binding sites



genomically informed starting place for
globally investigating detailed regulation

Main features of VISTA

- Clear , configurable output
- Ability to visualize several global alignments on the same scale
- Alignments up to several megabases
- Working with finished and draft sequences
- Available source code and WEB site

Related publications

O N E

Loots GG, Locksley RM, Blankespoor CM, Wang ZE, Miller W, Rubin EM, Frazer KA. Identification of a coordinate regulator of interleukins 4, 13, and 5 by cross-species sequence comparisons.

Science. 2000 Apr 7;288(5463):136-40.

T W O

Pennacchio LA, Olivier M, Hubacek JA, Cohen JC, Cox DR, Fruchart JC, Krauss RM, Rubin EM. An apolipoprotein influencing triglycerides in humans and mice revealed by comparative sequencing.

Science. 2001 Oct 5;294(5540):169-73.

What if you don't have a sequence of other species for the region of your interest?

Are there publicly available comparative genomics data?

Large scale VISTA applications:

Cardiovascular comparative genomics database

<http://pga.lbl.gov>

Godzilla - comparing the human and mouse genome

<http://pipeline.lbl.gov>



Godzilla - automatic computational system for comparative analysis of genomes

<http://pipeline.lbl.gov> <http://www-gsd.lbl.gov/vista>

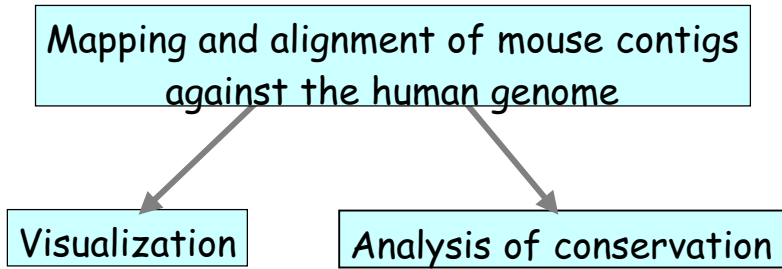
DATA

Base Human Genome - Golden Path Assembly

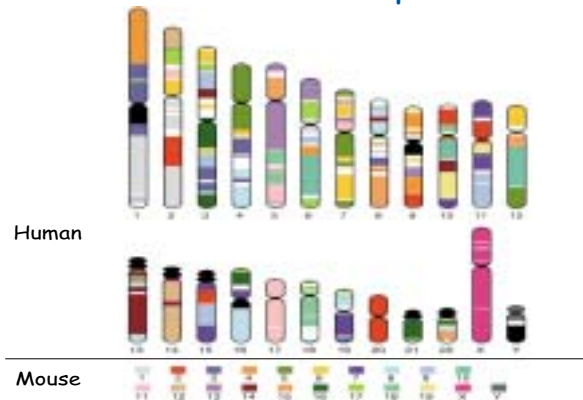
Mouse assemblies:

Arachne	October 2001
Phusion	November 2001
MGSC v3	April 2002

Main modules of the system



Chromosome Comparison



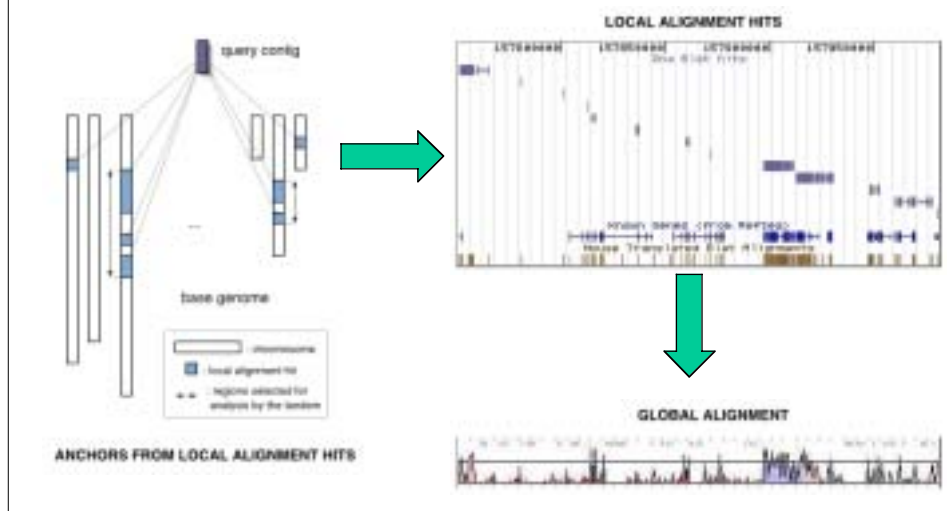
Base pair alignment

```
247 GGTGAGGTCGAGGACCCTGCA CCGAGCTGTATGGAGGGCA AGAGC
    |:  ||  ||||:  |||  -:||  |||  |:|  |||---|||
368 GAGTCGGGGGAGGGGGCTGCTGTTGGCTCTGGACAGCTTGCATTGAGAGG
```

Tandem Local/Global Alignment Approach

Sequence fragment **anchoring** (DNA and/or translated BLAT)

Multi-step verification of potential regions using global alignment (AVID)

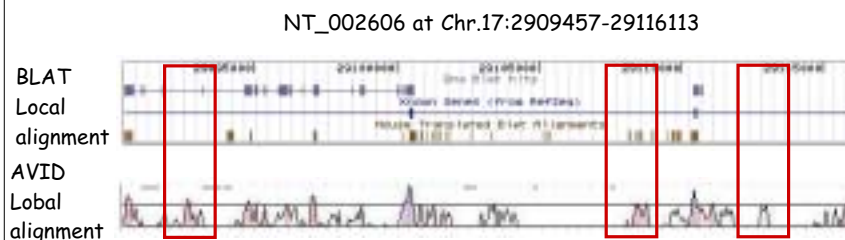


Advantage of the tandem approach:

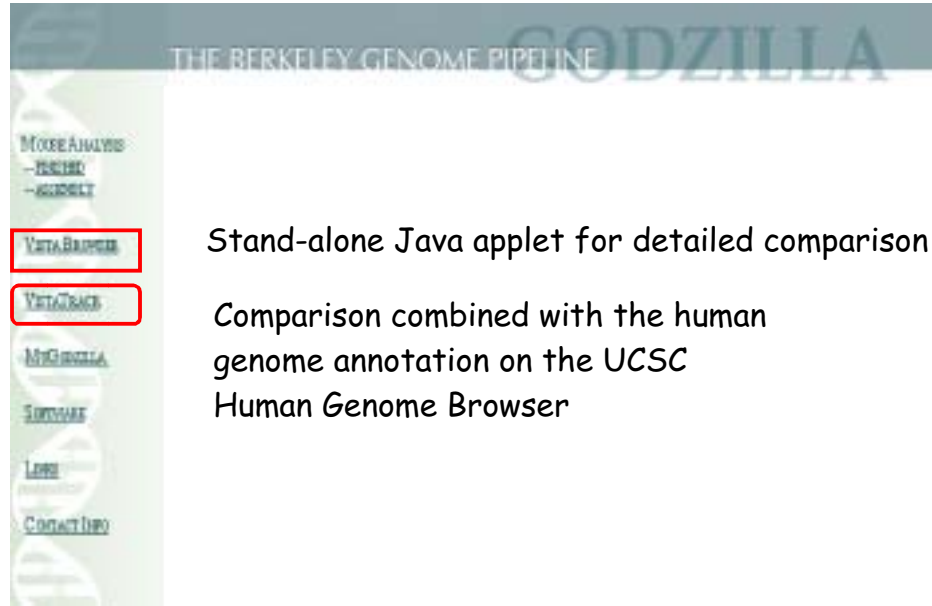
better sensitivity/specificity trade-off

fill-in effect

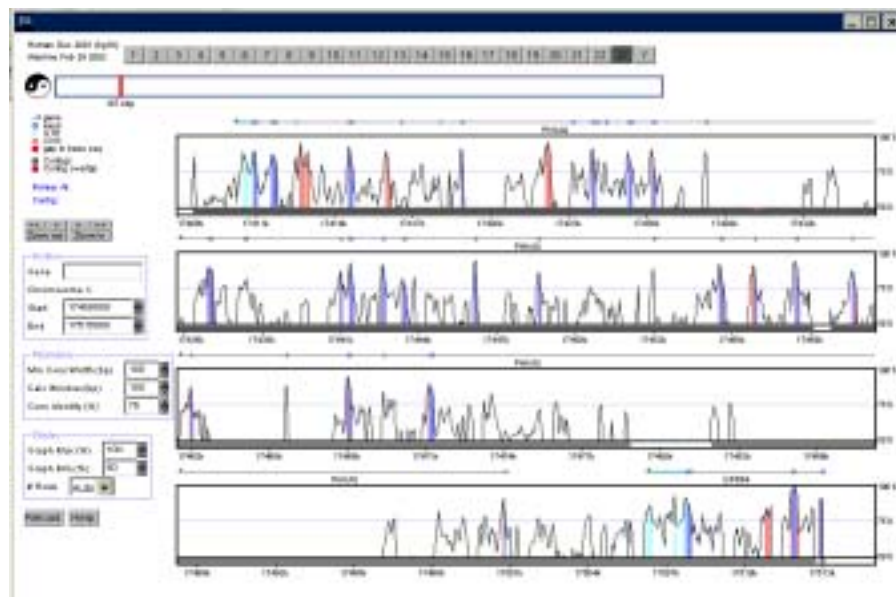
scoring longer alignments



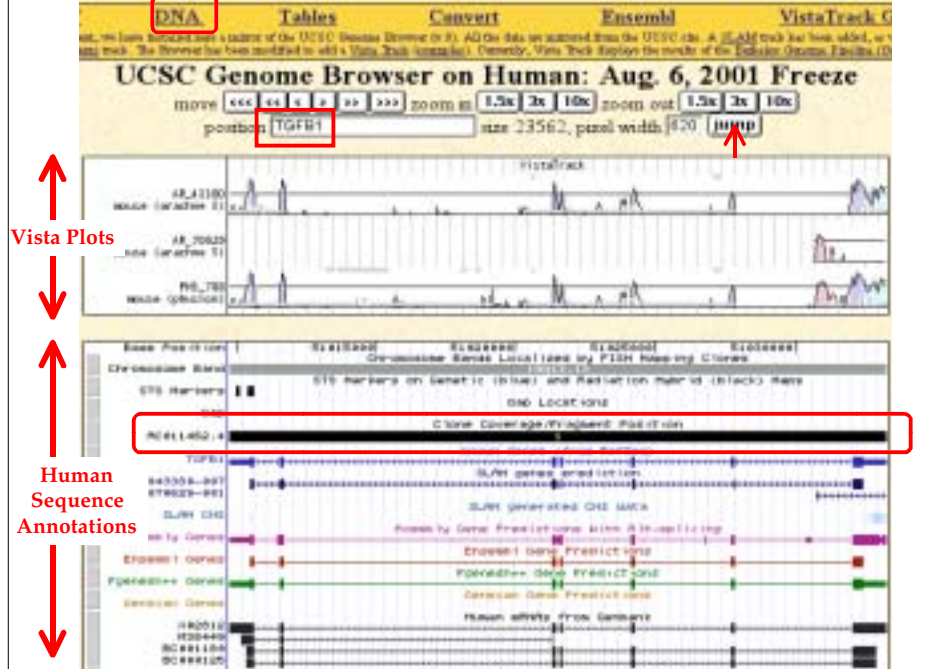
Visualization - VistaBrowser & VistaTrack



VistaBrowser



VistaTrack



<http://pipeline.lbl.gov/>



MyGodzilla Tool



Submit a DNA sequence of ANY organism...
... or submit a whole chromosome and analyze another Genome



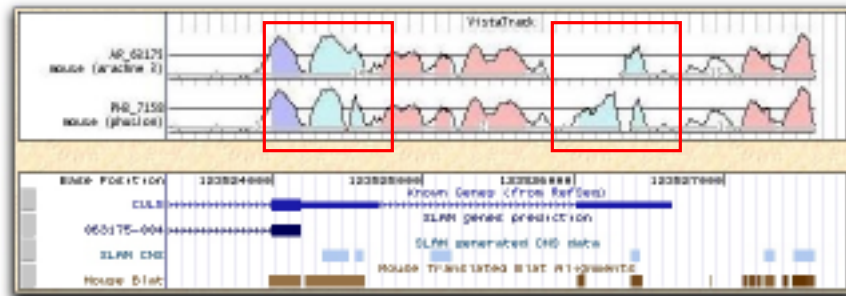
Your request is submitted and the results will be ready at this link: <http://pipeline.lbl.gov/cache24h/22599550.8711>

to receive an automatic notification type your email here:

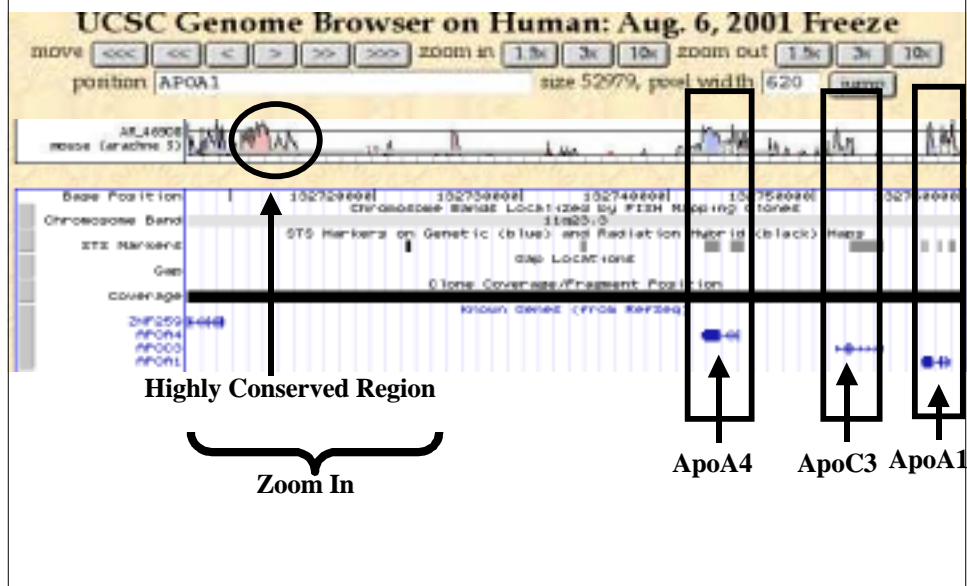
Examples of Results

- Understanding the structure of conservation
- Identification of putative functional sites
- Discovery of new genes
- Detection of contamination and misassemblies

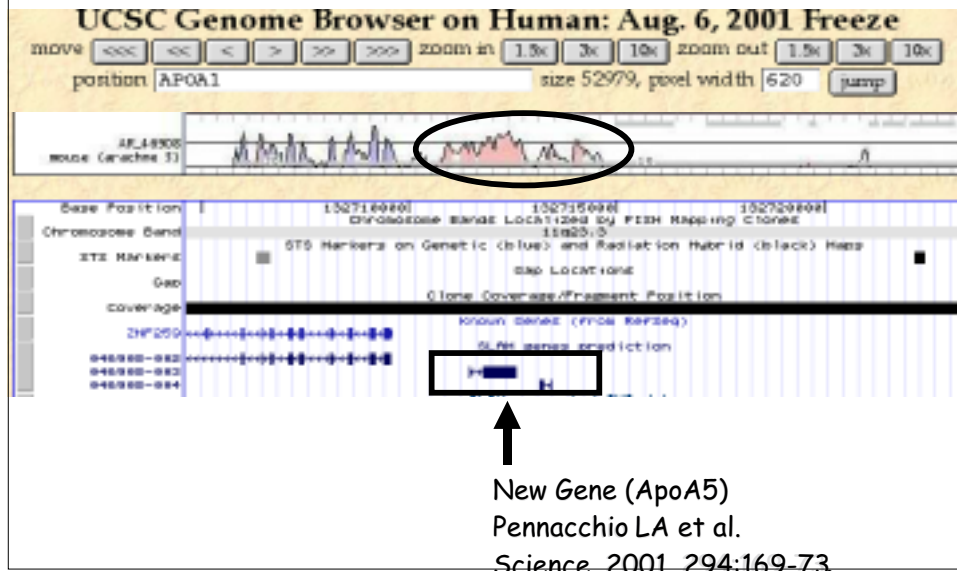
Two assemblies are better than one



Identification of a New Apo Gene on Human 11q23
Godzilla

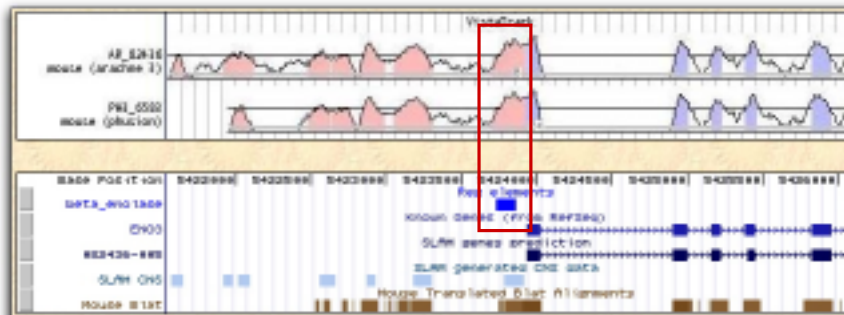


Identification of a New Apo Gene on Human 11q23 Godzilla



Finding regulatory regions

Muscle Specific Regulatory Region: human
beta enolase intronic enhancer



Comparative analysis of genomic intervals containing important cardiovascular genes

<http://pga.lbl.gov>



<http://pga.lbl.gov/cvcgd.html>



Example of CVCGD entry

Solute carrier family 22, organic cation transporter member 4 (SLC22A4, OCTN1) - RefSeq

Category: Amino acid transporter

Gene ID in the OMIM database: 604130

Human map location: 5q31

GenBank accession number for human cDNA: [U01001](#)

Mouse map location: 11

GenBank accession number for mouse cDNA: [U01002](#)

GenBank accession number for human cDNA: [U01003](#)

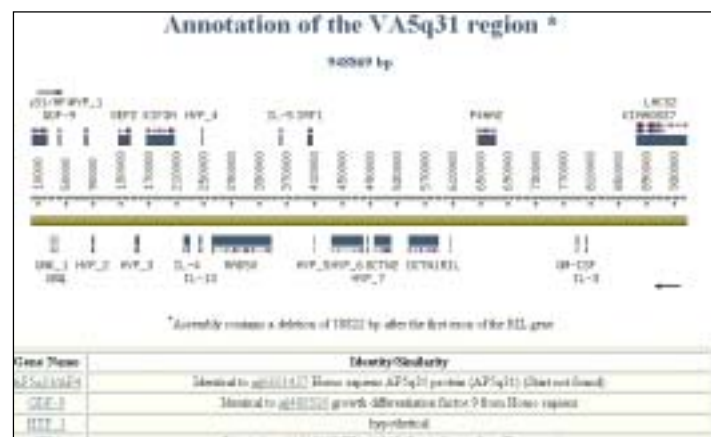
GenBank accession number for mouse cDNA: [U01004](#)

Human protein alignment: [Whole sequence](#) | [1-100000](#) | [100001-200000](#) | [200001-300000](#) | [300001-400000](#) | [400001-500000](#) | [500001-600000](#) | [600001-700000](#) | [700001-800000](#) | [800001-900000](#) | [900001-1000000](#) (see important note below) | [Download](#)

List of transmembrane regions

Note: If your browser hangs or crashes on the alignment page you can try [this link](#) instead.

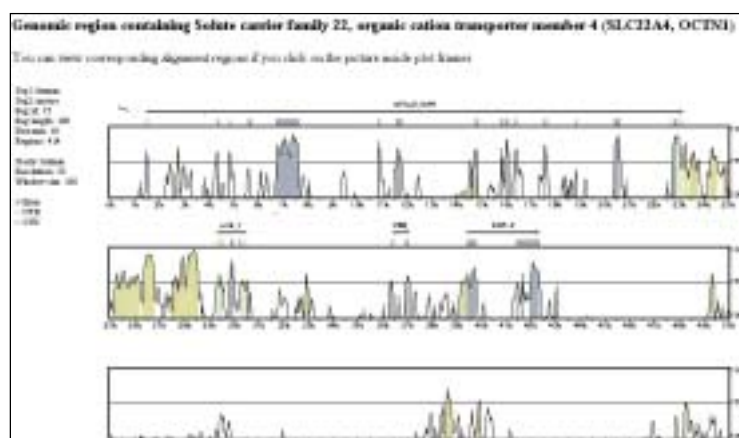
Short annotation of the region



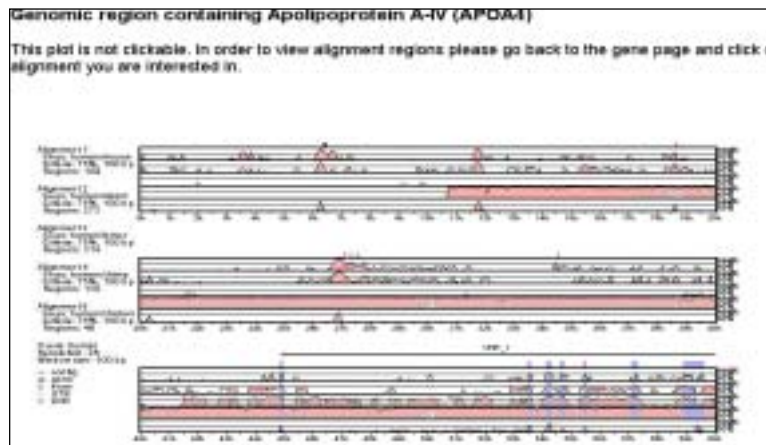
Detailed annotation in AceDB format



VISTA plot of the region



multiVISTA plot of the region



Alignment

Genomic region containing Solute carrier family 22, organic cation transporter membe

```
seq1 = human
seq2 = mouse

          6990      7000      7010      7020      7030      7040
seq1  CAGAGTGACAGCACACACAGAGGAGAACTGTAGCCCAAAABACAACCCAAAAGGCTGAG
      ||||| ||||| || || ||||| ||||| ||||| ||||| ||||| ||||| |||||
seq2  CAGAGCGGACATACCACCTCAGAGGAGAACTGTGCGCAAAABACAACCCAAAAGGCTGAG
      ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||

          8130      8140      8150      8160      8170      8180
seq1  AAGGCGAGCTGCTGAGAGGCGCTCGTGAGGCGCTGAGAGTAGAAGTGAAGACCCCTGTAGAC
      || ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
seq2  AAGTCAGCTGCTGAGAGGCGCTCGTGAGGCGCTGAGAGTAGAAGTGAAGACCCCTGTAGAC
      ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||

          8190      8200      8210      8220      8230      8240
seq1  TTGCGTAGCAGCATGCGCTCCAGCAGGACACAAAGCAGCCAGCCAGAGCGCTCAAGGAGACCC
      ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
seq2  ATGCGCTGCAGCATGCGCTCCAGCAGGACACAAAGCAGCCAGCCAGAGCGCTCAAGGAGACCC
      ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||

          7110      7120      7130      7140      7150      7160
seq1  TTGCGTAGCAGCATGCGCTCCAGCAGGACACAAAGCAGCCAGCCAGAGCGCTCAAGGAGACCC
      ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
seq2  ATGCGCTGCAGCATGCGCTCCAGCAGGACACAAAGCAGCCAGCCAGAGCGCTCAAGGAGACCC
      ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||

          8250      8260      8270      8280      8290      8300
```

Conserved regions

Genomic region containing Solute carrier family 22, organic cation transporter member 4 (SLC22A4, OCTN1)

Criteria: 70% identity over 100 bp

Conserved Regions - 3000 (3000)

1468	(1880)	to	1518	(1828)	=	47bp	at	88.1%	cons
2446	(2945)	to	2517	(2191)	=	175bp	at	89.4%	cons
4516	(4851)	to	4579	(4888)	=	53bp	at	100.0%	cons
4818	(4828)	to	4833	(4171)	=	55bp	at	87.4%	cons
6717	(7860)	to	7634	(5777)	=	238bp	at	87.0%	cons
12818	(12748)	to	12827	(12817)	=	9bp	at	81.0%	cons
11253	(12527)	to	11793	(12071)	=	235bp	at	81.0%	cons
14808	(18706)	to	14622	(18623)	=	118bp	at	76.5%	cons
14871	(18888)	to	14782	(18603)	=	118bp	at	74.5%	cons
18784	(11693)	to	18878	(11609)	=	235bp	at	88.5%	cons
18797	(17124)	to	18848	(17088)	=	60bp	at	81.0%	cons
15975	(17923)	to	16111	(17816)	=	127bp	at	88.5%	cons
14348	(18848)	to	14434	(18116)	=	118bp	at	81.7%	cons
18417	(18517)	to	18525	(18217)	=	108bp	at	78.1%	cons
17554	(18816)	to	17647	(18007)	=	54bp	at	87.1%	cons

Publications on our tools:

- I. Dubchak, M. Brudno, L.S. Pachter, G.G. Loots, C. Mayor, E. M. Rubin, K. A. Frazer. (2000) Active conservation of noncoding sequences revealed by 3-way species comparisons. *Genome Research*, 10: 1304-1306.
- C. Mayor, M. Brudno, J. R. Schwartz, A. Poliakov, E. M. Rubin, K. A. Frazer, Lior S. Pachter, I. Dubchak. (2000) VISTA: Visualizing global DNA sequence alignments of arbitrary length. *Bioinformatics*, 16: 1046-1047.
- G. G. Loots, I. Ovcharenko, L. Pachter, I. Dubchak and E. M. Rubin. (2002) Comparative sequence-based approach to high-throughput discovery of functional regulatory elements. *Genome Res.*, 12:832-839
- I. Dubchak, L. Pachter. (2002) The computational challenges of applying comparative-based computational methods to whole genomes. *Briefings in Bioinformatics*, 3, 18.

Related sites

- The Human Genome Browser & BLAT program
<http://genome.ucsc.edu/>
- ENSEMBLE Project (Sanger Center) <http://www.ensembl.org/>
- AVID alignment program
<http://baboon.math.berkeley.edu/~syntenic/avid.html>
- SLAM comparative gene prediction program
<http://bio.math.berkeley.edu/slam/mouse/>
- PSU group's MHC Human-Mouse comparison results
<http://bio.cse.psu.edu/mousegroup/MHC/>
- PSU Pipmaker program <http://bio.cse.psu.edu/pipmaker/>

Summary

Suite of comparative genomics tools **VISTA**
<http://www-gsd.lbl.gov>

Godzilla comparing the human and mouse genome
<http://pipeline.lbl.gov>

Cardiovascular comparative genomics database
<http://pga.lbl.gov>

Questions? Write to vista@lbl.gov

Towards Better VISTAs

Information
from a Single
Sequence
Alone



Multi-Organism
High Quality
Sequences



Thanks

Biology

Kelly Frazer
Gaby Loots
Len Pennacchio

Eddy Rubin

Bioinformatics

Michael Brudno
Olivier Couronne
Brian Klock
Chris Mayor
Ivan Ovcharenko
Alexander Poliakov
Jody Schwartz
Lior Pachter (UCB)

Funding - Programs for Genomic Applications (PGA) by NHLBI